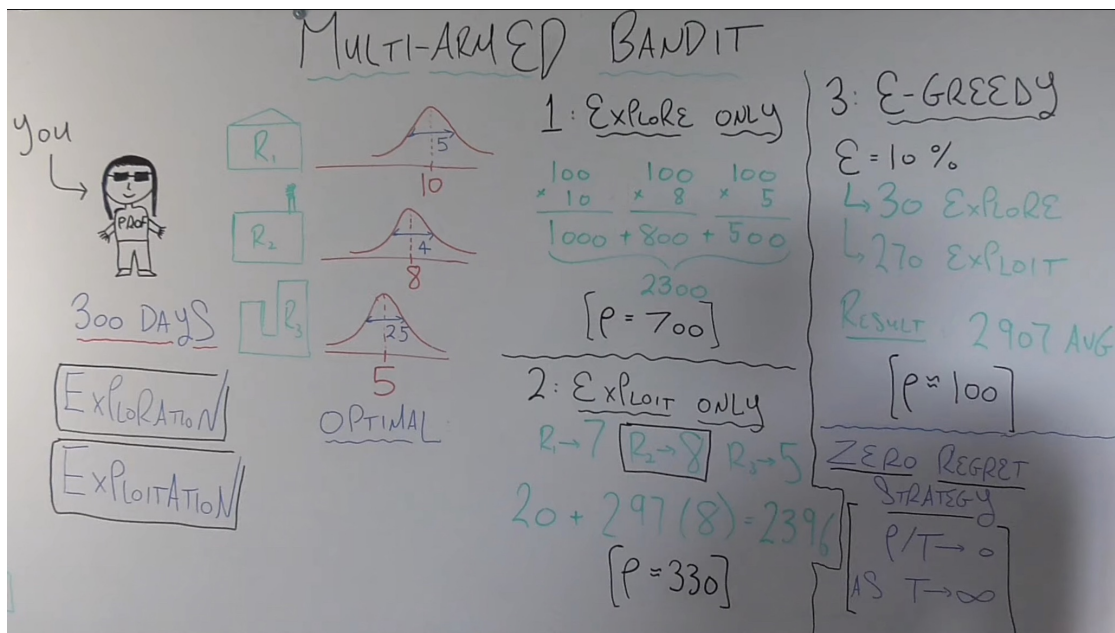




Fundamentals of Reinforcement Learning



$P \rightarrow$ regret \rightarrow Diff. b/w max happiness available - happiness from your strategy

'sho'

$$q^*(a) = E[R_t | A_t = a]$$

↓
value of action

↘ action selected on time step t

$Q_t(a) \rightarrow$ estimated value of action a at time step t

$$Q_{n+1} = Q_n + \alpha [R_n - Q_n]$$



Increment Update rule

New estimate = Old estimate + stepsize

[Target - Old Estimate]

Q_n for nonstationary Problem

$$Q_{n+1} = (1 - \alpha)^n Q_1 + \sum_{i=1}^n \alpha (1 - \alpha)^{n-i} R_i$$

$\Rightarrow A_t \rightarrow$ action at time t

Ways of balancing Exploration & Exploitation

I Epsilon greedy

$$A_t \leftarrow \begin{cases} \operatorname{argmax}_a Q_t(a) & \text{with probability } 1 - \epsilon \\ a \sim \text{Uniform}(\{a_1, \dots, a_k\}) & \text{with probability } \epsilon \end{cases}$$

exploitation

exploration (random)

II Optimistic initial value

III

Upper - Confidence Bound (UCB) Action Selection

$$A_t = \operatorname{argmax} \left[Q_t(a) + c \sqrt{\frac{\ln t}{N_t(a)}} \right]$$

Annotations:

- $Q_t(a)$: exploit
- $c \sqrt{\frac{\ln t}{N_t(a)}}$: uncertainty term (circled), explore

$$c \sqrt{\frac{\ln t}{N_t(a)}} \rightarrow c \sqrt{\frac{\ln \text{timesteps}}{\text{times action a taken}}}$$

Examples:

- $c \sqrt{\frac{\ln 10000}{5000}} \rightarrow 0.043c$
- $c \sqrt{\frac{\ln 10000}{100}} \rightarrow 0.303c$

↑
less time → more uncertainty

IV Gradient Bandit

X X X X

Markov Decision Process

$$p(s', r | s, a) = \Pr\{S_t = s', R_t = r \mid S_{t-1} = s, A_{t-1} = a\}$$

$S_t \in S$ \rightarrow set of possible states
 $A_t \in A(s_t)$ \rightarrow set of possible actions at that state

$$\sum_{s' \in S} \sum_{r \in R} p(s', r | s, a) = 1 \text{ for all } s \in S, a \in A(s)$$

$$r(s, a) = E[R_t \mid S_{t-1} = s, A_{t-1} = a] =$$

$$r(s, a, s') = E[R_t \mid S_{t-1} = s, A_{t-1} = a, S_t = s']$$

$0 \leq \gamma \leq 1 \implies$ Discount rate

Return at timestep t (episodic tasks)

$$G_t = R_{t+1} + R_{t+2} + \dots + R_T \quad \text{terminal state}$$

Continuous Tasks

$$G_t = R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3} + \dots + \gamma^{k-1} R_{t+k} + \dots$$
$$= \sum_{k=0}^{\infty} \gamma^k R_{t+k+1}$$

$$G_t = R_{t+1} + \gamma G_{t+1} \quad \text{recursive definition}$$

If γ approaches 1 \longrightarrow long term
(take future rewards into more consideration)

$$G_t = \sum_{k=0}^{\infty} \gamma^k = \frac{1}{1-\gamma}$$

$$G_t = \sum_{k=0}^{\infty} \gamma^k R_{t+k+1} \leq \sum_{k=0}^{\infty} \gamma^k R_{\max} = R_{\max} \sum_{k=0}^{\infty} \gamma^k$$

$$= R_{\max} \times \frac{1}{1-\gamma}$$

constant

(it's not changing)
(every action gets R_{\max})

Policies

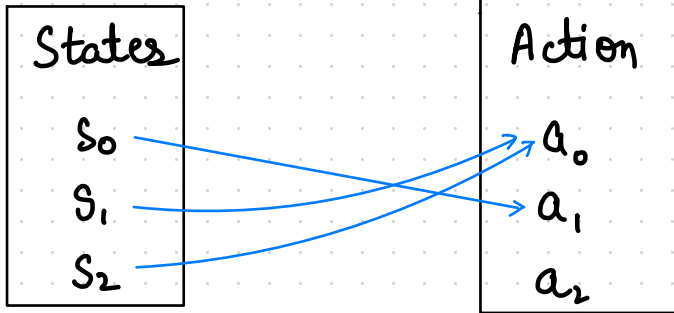
→ specifies how an agent behaves

Policy is a distribution over actions for each possible state.

→ A policy maps each state to a single action.

$$\pi(s) = a$$

→ represents the action selected in state 's' by policy ' π '.



Deterministic
policy
notation

In general policy assign probabilities to each action in each state.

$\pi(a|s)$ → probability of selecting action a in a state s .

$$\sum_{a \in A(s)} \pi(a|s) = 1 ; \pi(a|s) \geq 0$$

$a \in A(s)$

Stochastic policy notation

Value Functions

I State Value Functions

$$V_{\pi}(s) \doteq E_{\pi}[G_t | S_t = s] = E_{\pi}\left[\sum_{k=0}^{\infty} \gamma^k R_{t+k+1} | S_t = s\right]$$

$R_{t+1} + \gamma R_{t+2} + \gamma^2 R_{t+3}$
↑

a state value function is the future award an agent can expect to receive starting from a particular state.


- (expected return from a given state)

II Action - value Function

$$G_t = R_t + \gamma G_{t+1}$$

$$Q_{\pi}(s, a) \doteq E_{\pi}[G_t | S_t = s, A_t = a]$$

an action value describes what happens when the agent first selects a particular action.

- (the action value of a state is the expected return if the agent selects action A and then follows policy ' π ')


State value functions represent the expected return from a given state under a specific policy

Action-value functions represent the expected return from a given state after taking a specific action, later following a specific policy



In reinforcement learning, "value functions" are a fundamental concept used to estimate and evaluate the desirability of different states and state-action pairs within an environment. Value functions play a crucial role in guiding an agent's decision-making process and learning how to maximize cumulative rewards. There are two main types of value functions:



1. State-Value Function (V-function, V):

- The state-value function, denoted as $V(s)$, estimates the expected return or cumulative reward an agent can achieve starting from a specific state while following a particular policy.
- $V(s)$ quantifies how good it is to be in a particular state under a given policy. It represents the long-term expected reward the agent can expect if it starts in state s and continues to interact with the environment according to the policy.
- The formula for the state-value function is typically expressed as:

$$V(s) = E \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s \right]$$

where E denotes the expectation, γ is a discount factor that gives less weight to distant rewards, R_{t+1} is the reward received at time step $t + 1$, and S_0 is the initial state.

2. Action-Value Function (Q-function, Q):

- The action-value function, denoted as $Q(s, a)$, estimates the expected return or cumulative reward an agent can achieve when it starts from a specific state, takes a particular action, and then continues to interact with the environment following a specific policy.
- $Q(s, a)$ quantifies the desirability of taking a particular action, a , in a specific state, s , under the given policy.
- The formula for the action-value function is often expressed as:

$$Q(s, a) = E \left[\sum_{t=0}^{\infty} \gamma^t R_{t+1} \mid S_0 = s, A_0 = a \right]$$

where A_0 is the initial action taken.

Regenerate

Value functions provide the agent with a way to assess and compare different states or state-action pairs in terms of their expected long-term rewards. The goal of reinforcement learning is to find the optimal policy (the best strategy) that leads to the highest values in states, and, in turn, maximizes cumulative rewards. Various algorithms, such as dynamic programming, Monte Carlo methods, and temporal difference learning, are used to estimate and update value functions to achieve this goal.

Bellman Equation

State-value Bellman equation

Bellman equation for a state value function defines a relationship between the value of a state and the value of his possible successor states.

$$\begin{aligned} V_{\pi}(s) &\doteq E_{\pi}[G_t | S_t = s] \\ &= E_{\pi}[R_{t+1} + \gamma G_{t+1} | S_t = s] \\ &= \sum_a \pi(a, s) \sum_{s'} \sum_{\mathcal{R}} p(s', \mathcal{R} | s, a) [\mathcal{R} + \gamma E_{\pi}[G_{t+1} | S_{t+1} = s']] \\ &= \sum_a \pi(a, s) \sum_{s'} \sum_{\mathcal{R}} p(s', \mathcal{R} | s, a) [\mathcal{R} + \gamma V_{\pi}(s')] \end{aligned}$$

Action-value Bellman equation

$$q_{\pi}(s, a) = E_{\pi}[G_t | S_t = s, A_t = a]$$

action is already fixed

$$\begin{aligned} &= \sum_{s'} \sum_{\mathcal{R}} p(s', \mathcal{R} | s, a) [\mathcal{R} + \gamma E_{\pi}[G_{t+1} | S_{t+1} = s']] \\ &= \sum_{s'} \sum_{\mathcal{R}} p(s', \mathcal{R} | s, a) \left[\mathcal{R} + \gamma \sum_{a'} \pi(a' | s') E_{\pi}[G_{t+1} | S_{t+1} = s', A_{t+1} = a'] \right] \\ &= \sum_{s'} \sum_{\mathcal{R}} p(s', \mathcal{R} | s, a) \left[\mathcal{R} + \gamma \sum_{a'} \pi(a' | s') q_{\pi}(s', a') \right] \end{aligned}$$

Optimal Policy

An optimal policy π_* is a policy which is as good as or better than all the other policies.

$\Rightarrow \pi_1 \geq \pi_2$ if & only if $V_{\pi_1}(s) \geq V_{\pi_2}(s)$ for all $s \in S$.

Optimal value function

$$V_* \quad V_{\pi_*}(s) \doteq \sum_{\pi_*} [G_t | S_t = s] = \underline{\max_{\pi} V_{\pi}(s)} \text{ for all } s \in S$$

$$Q_* \quad Q_{\pi_*}(s, a) = \max_{\pi} Q_{\pi}(s, a) \text{ for all } s \in S \text{ and } a \in A$$

$$V_{\pi}(s) = \sum_a \pi(a, s) \sum_{s'} \sum_a p(s', r | s, a) [r + \gamma V_{\pi}(s')] \quad (\text{bellman eq.})$$

$$V_*(s) = \sum_a \pi_*(a, s) \sum_{s'} \sum_a p(s', r | s, a) [r + \gamma V_*(s')]$$

$$V_*(s) = \max_a \sum_{s'} \sum_a p(s', r | s, a) [r + \gamma V_*(s')]$$

Bellman Optimality Equation for V_*

$$q_*(s,a) = \sum_{s'} \sum_{\mathbf{r}} p(s', \mathbf{r} | s, a) \left[\mathbf{r} + \gamma \max_{a'} q_*(s', a') \right]$$

Bellman Optimality Equation for q_*

$$V_*(s) = \max_a \sum_{s'} \sum_a p(s', a | s, a) [r + \gamma V_*(s')]$$

$$\pi_*(s) = \operatorname{argmax}_a \sum_{s'} \sum_a p(s', a | s, a) [r + \gamma V_*(s')]$$

$$r + \gamma V_*(s')$$

$$\pi_*(s) = \operatorname{argmax}_a q_*(s, a)$$

DYNAMIC PROGRAMMING

(week 4)

$$[x + \gamma \overline{V_{\pi}(s')}]^D$$

$$\boxed{-1}$$

$$x = \frac{3}{4}(1-x) + \frac{1}{4}(-2)$$

$$-1 + \gamma V_{\pi}(s')$$

$$-1 + [-x] = \frac{1}{4} \quad \begin{matrix} [-1 + (-2)] \\ x = 0.25(-2) \\ x = \end{matrix}$$

$$x = \frac{1}{4} [x + \gamma V(s')]$$

$$= \frac{1}{4} [-1 + x]$$

$$\approx 0.25 + 0.25x$$